

A visualization of the cosmic web, showing a complex network of blue filaments and red nodes against a dark background.

Spectroscopic Databases and Manifold Learning for Surveys of the 2020s

Dan Masters

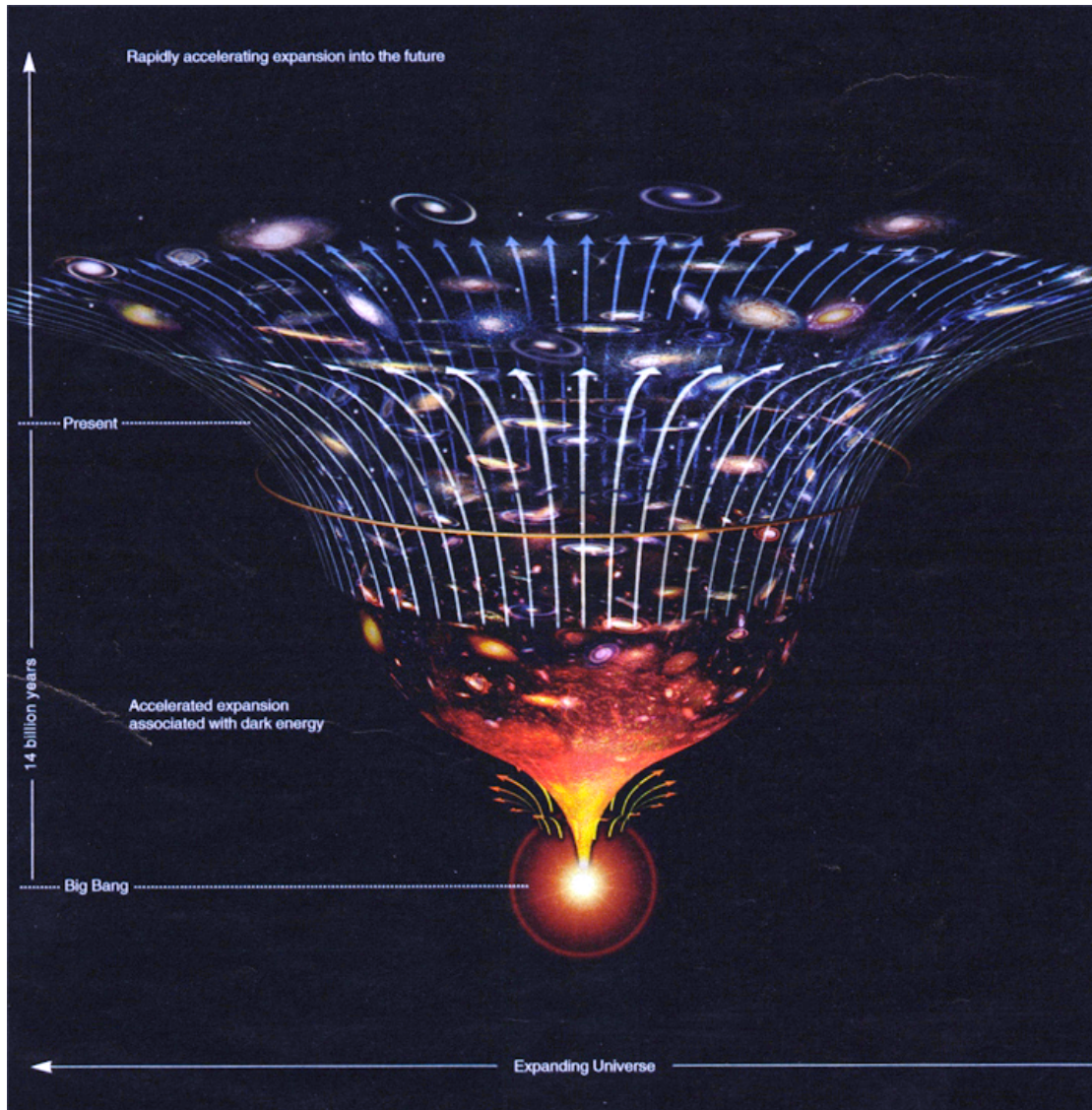
Jet Propulsion Laboratory, California Institute of Technology

December 6, 2018



Jet Propulsion Laboratory
California Institute of Technology

© 2018 California Institute of Technology
Government sponsorship acknowledged.



Breakthrough from the 1990s:
Accelerating cosmic expansion

2011 Nobel Prize in Physics



Λ ?

The Redshift Measurement Problem

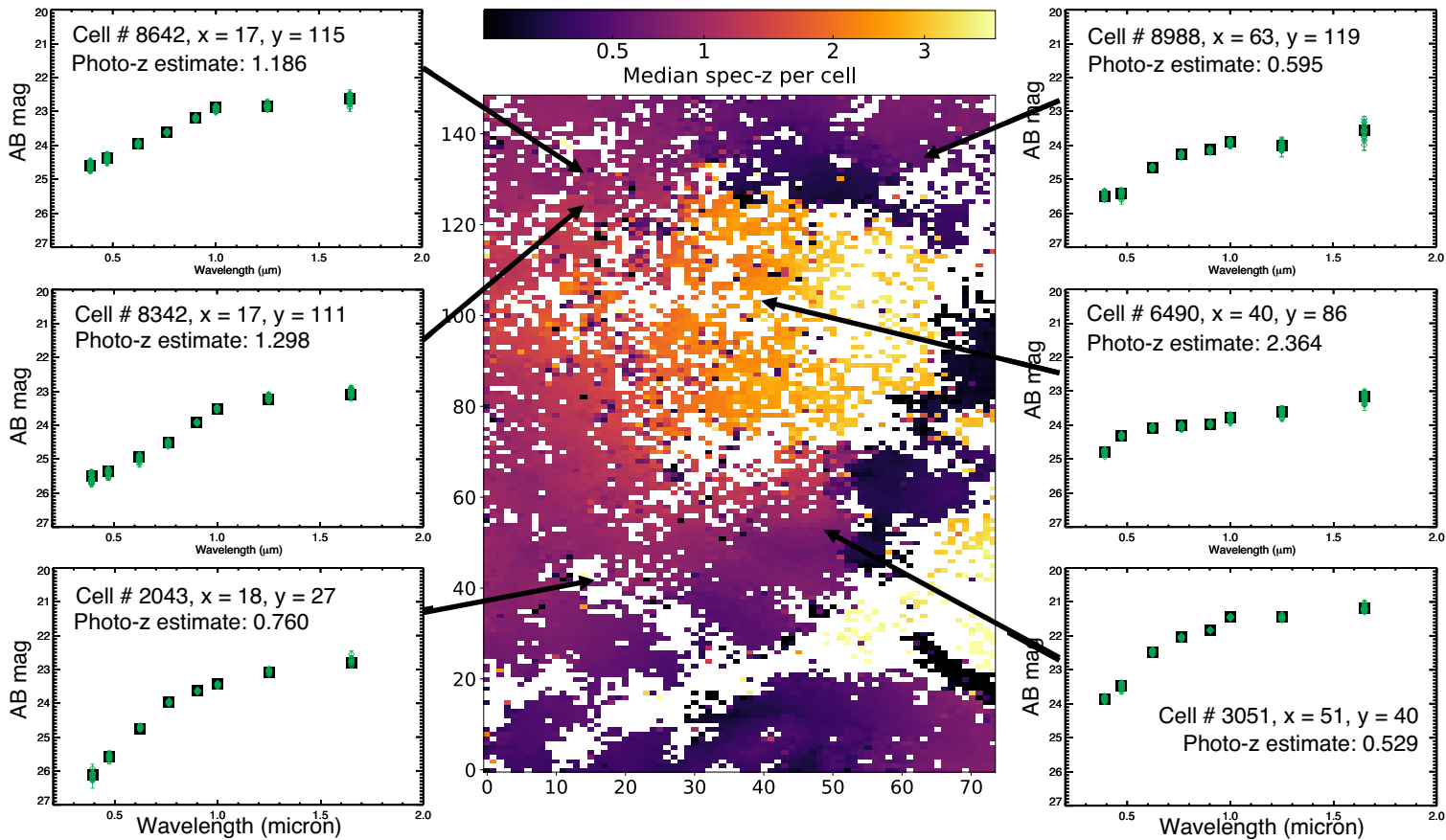
- Billions of galaxies will be imaged by the Stage IV cosmology surveys (LSST, Euclid, WFIRST)
- Only possible to get spectroscopic redshifts for a small fraction
- Weak lensing cosmology requires that the redshift distributions of galaxies in ~ 10 - 20 redshifts bins be known with **high accuracy**

→ **Photometric redshifts will necessarily be crucial for weak lensing cosmology missions**

Manifold learning / nonlinear dimensionality reduction (NLDR)

- Group of techniques to characterize / explore high-dimensional data and correlations in high dimensions
- Common ones includes the self-organizing map (SOM), t-SNE, local linear embedding (LLE), and UMap
- Most project the high-D manifold down to a lower-D representation
- Whereas deep convolution networks try to learn a complex high-dimensional relationship between input data and output labels, NLDR just tries to unwrap the high-D data in an unsupervised way – no outputs

Self-organized map of galaxy colors to Euclid depth



Masters et al. 2015, ApJ, 813, 53

The galaxy color manifold

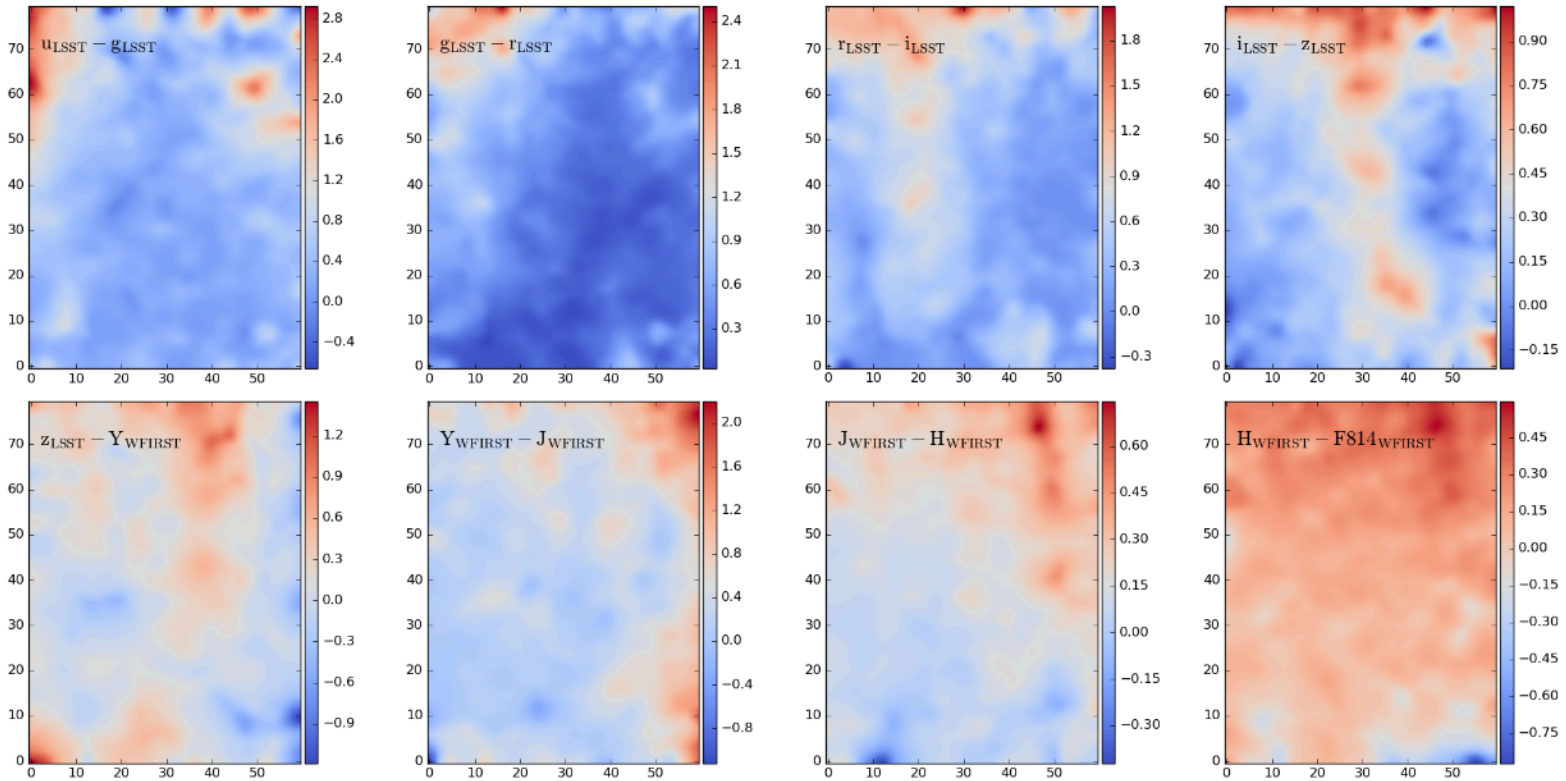
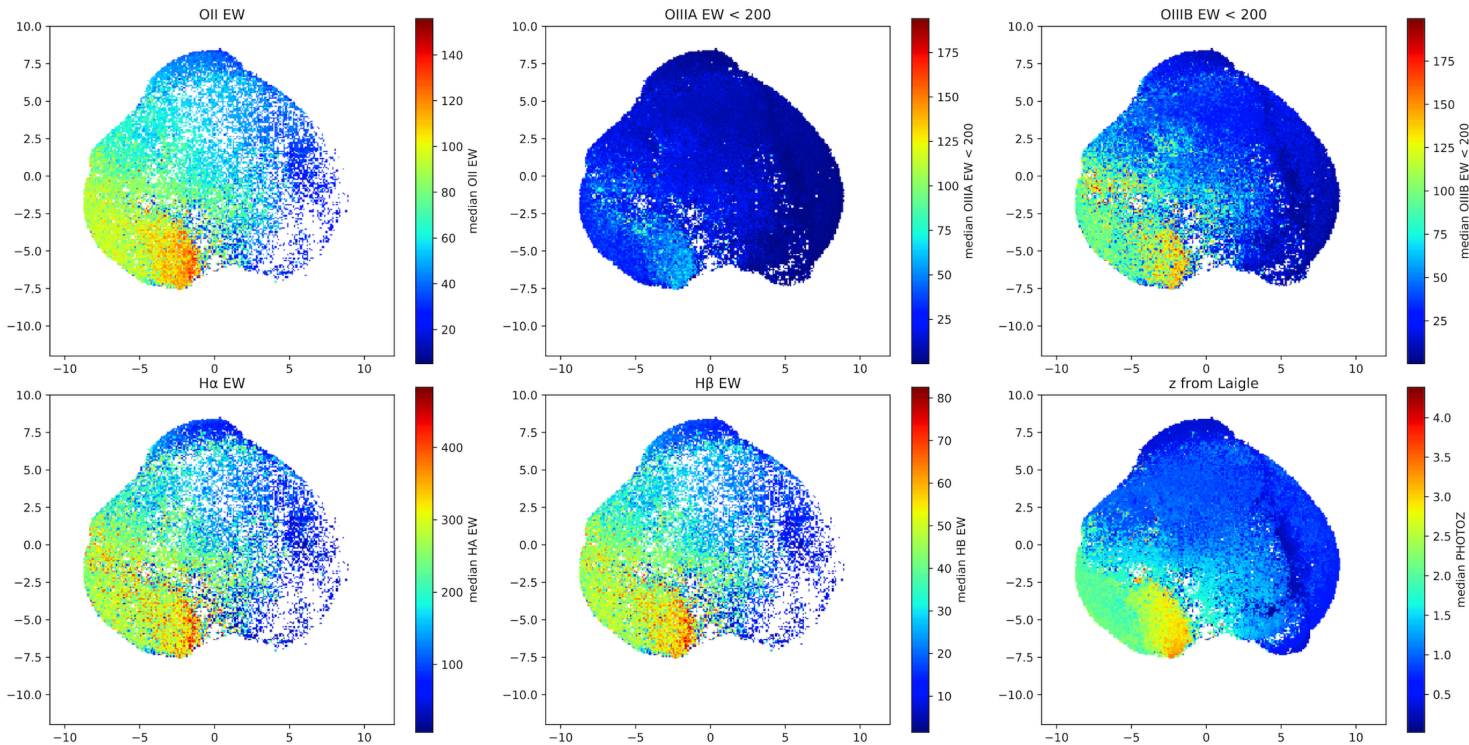


Figure 5. LSST and *WFIRST* colors of the trained SOM at each cell from top-left to bottom-right color-coded by: $u_{\text{LSST}} - g_{\text{LSST}}$, $g_{\text{LSST}} - r_{\text{LSST}}$, $r_{\text{LSST}} - i_{\text{LSST}}$, $i_{\text{LSST}} - z_{\text{LSST}}$, $z_{\text{LSST}} - Y_{\text{WFIRST}}$, $Y_{\text{WFIRST}} - J_{\text{WFIRST}}$, $J_{\text{WFIRST}} - H_{\text{WFIRST}}$, and $H_{\text{WFIRST}} - F814_{\text{WFIRST}}$. SOM is selected to be a mesh of 80×60 cells. The axes are arbitrary and each position on the two dimensional map points to a position in the 8 dimensional color space.

Hemmati et al. 2018

Other techniques – UMap



C3R2 = Complete Calibration of the Color-Redshift Relation

Judith Cohen (Caltech) - PI of Caltech Keck C3R2 allocation

16 nights (DEIMOS + LRIS + MOSFIRE, [kicked off program in 2016A](#))

Daniel Stern (JPL) - PI of NASA Keck C3R2 allocation

10 nights (all DEIMOS; “Key Strategic Mission Support”)

Daniel Masters (JPL) – PI of NASA Keck C3R2 allocation 2018A/B. [Observed last night, and will again tonight](#)

10 nights (5 each LRIS/MOSFIRE; “Key Strategic Mission Support”)

Dave Sanders (IfA) - PI of Univ. of Hawaii Keck C3R2 allocation

6 nights (all DEIMOS) + H20

Bahram Mobasher (UC-Riverside) - PI of UC Keck C3R2 allocation

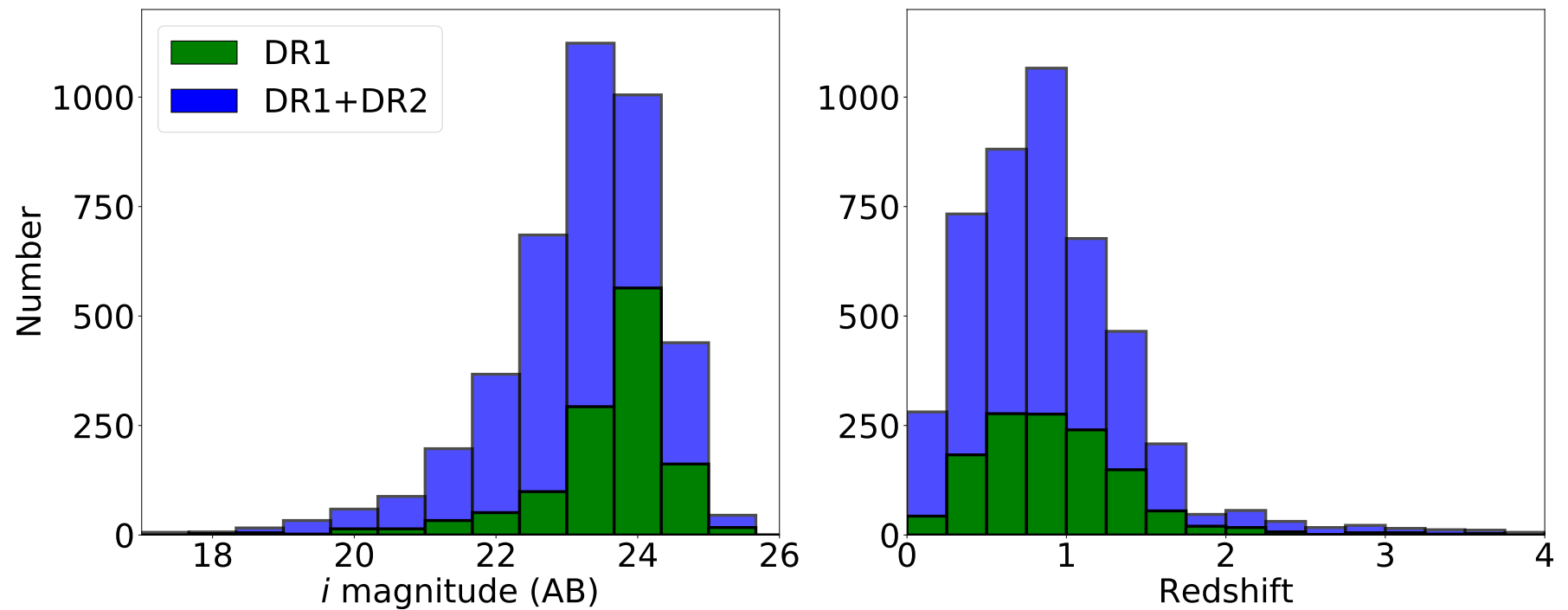
2.5 nights (all DEIMOS)

+ time allocations on VLT (PI F. Castander), MMT (PI D. Eisenstein), and GTC (PI C. Guitierrez)

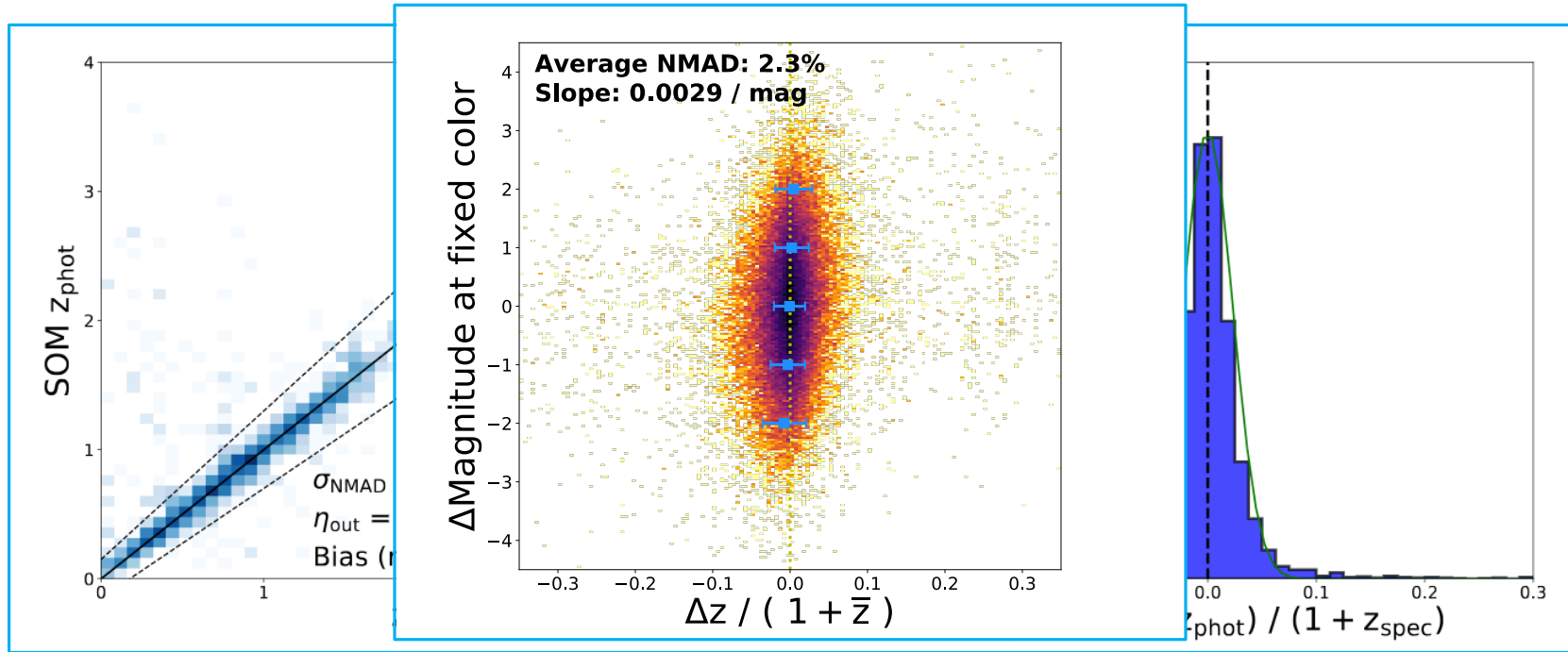
-Sample drawn from 6 fields totaling $\sim 6 \text{ deg}^2$

Additional Collaborators: Peter Capak, S. Adam Stanford, Nina Hernitschek, Francisco Castander, Sotiria Fotopoulou, Audrey Galametz, Iary Davidzon, Stephane Paltani, Jason Rhodes, Alessandro Rettura, Istvan Szapudi, and the Euclid Organization Unit – Photometric Redshifts (OU-PHZ) team

C3R2-Keck stats through DR2 (2016A-2017A)

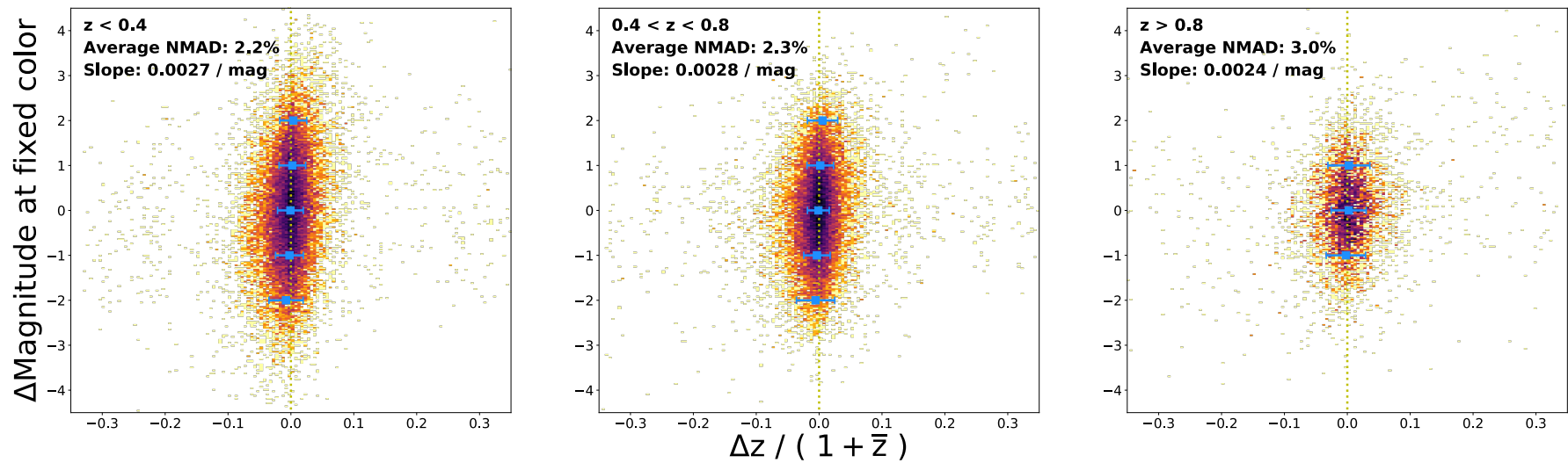


C3R2 – Results from SOM method

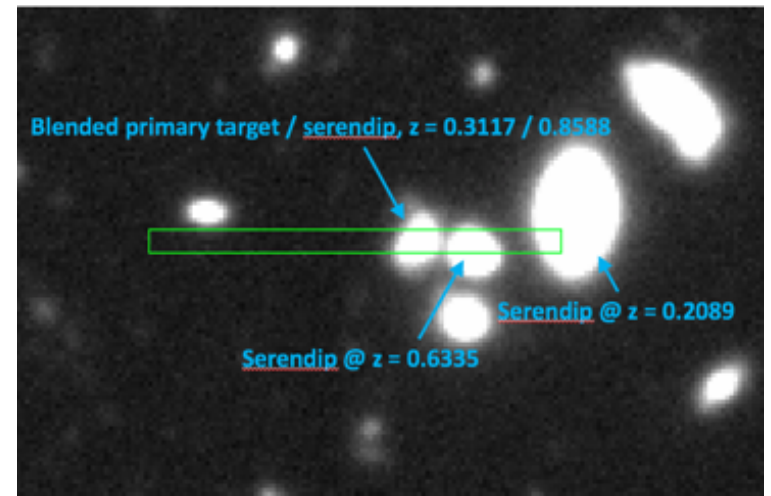
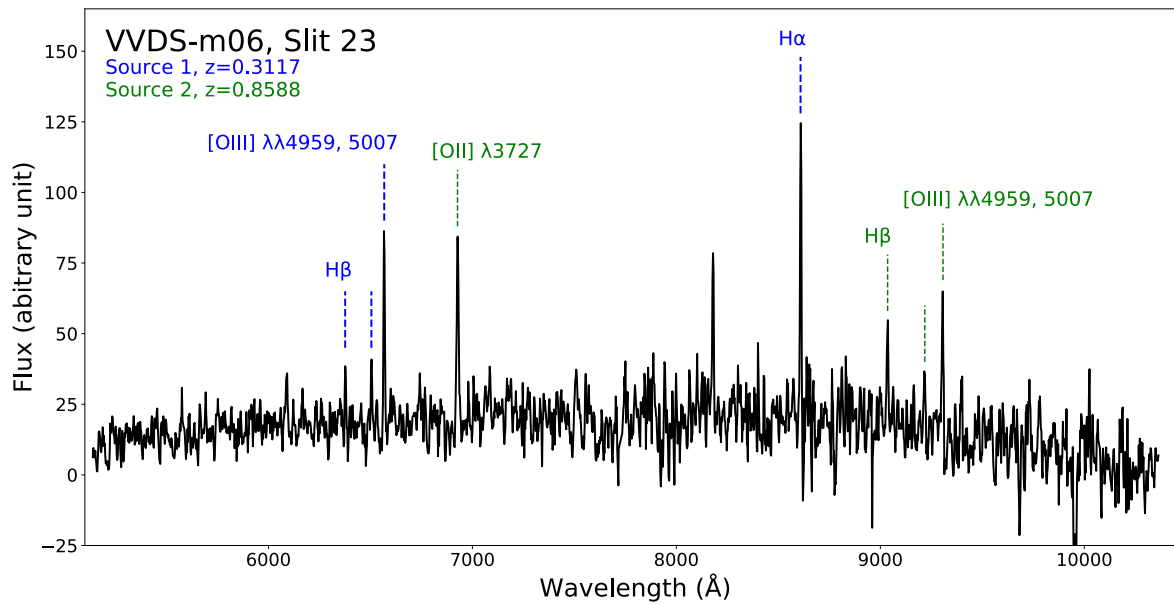


- Compare spectral energy distribution (SED) fitting performance to SOM position)
 - Method achieving unbiased performance
- Illustrates weak (and measurable) magnitude dependence of redshift on magnitude *at fixed color* in the LSST+Euclid color space

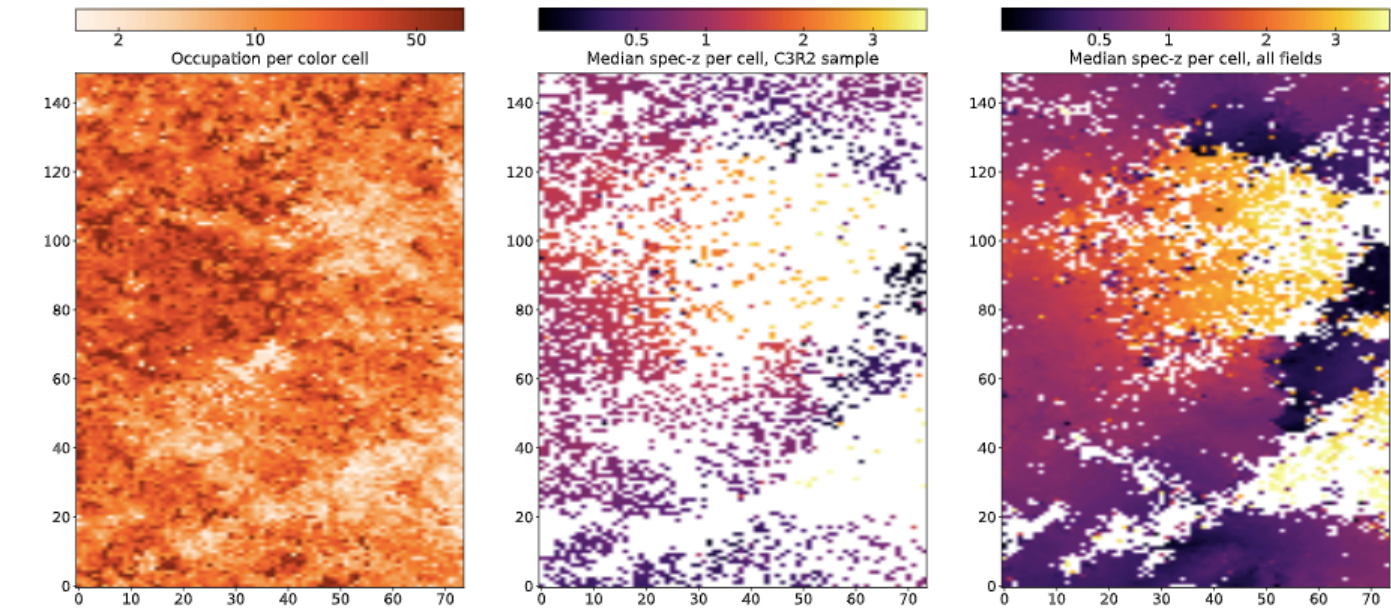
Remarkably stable relationship of $dmag/dz$ at fixed color



Ambiguous redshifts



Color coverage to Euclid depth



- Much of the galaxy color space explored to significant depth (>75% of cells; >85% of galaxies in cell with at least 1 specz, many cells with >>1 specz); some cells remain uncalibrated at present.
- C3R2-Keck alone covering >35% of the color space

Spectroscopic Databases – challenges

- We will have hundreds of thousands of deep galaxy spectra in the 2020s, growing continuously
- Need a database that can easily ingest new spectroscopy from disparate sources (e.g., from grisms)
- Extremely careful vetting of spec-zs necessary for cosmology
- Machine learning-based redshifts may prove critical

All large cosmology surveys would benefit from a single high-quality database of all deep spectroscopy!

A start for Euclid in Geneva

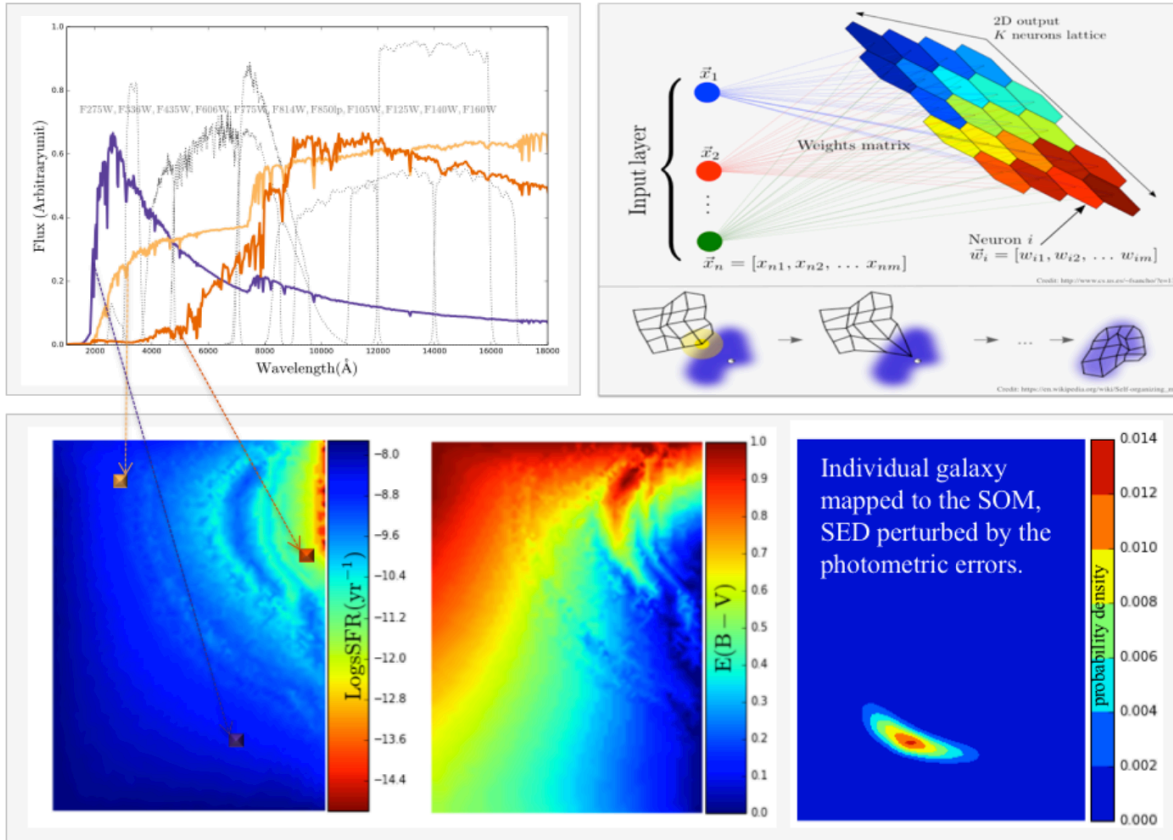
The database – Case example / Read-only



Galaxies are not unique

- The manifold of galaxy observables is finite, and large surveys like Euclid/WFIRST/LSST will measure essentially the same galaxy over and over
- We can measure the galaxy manifold really well with large surveys.
- Continuity constraints could then allow us to build a dynamic picture of galaxy growth
 - Individual galaxies can be thought of as moving along the manifold.
- What could we learn from this?

Models – can we match them to the data?



Measure the high-dimensional manifold.
Then what?

- We have a well-defined target for simulations
- What if we find (as is common) that the simulations produce unphysical galaxies, or can't produce certain real galaxies?
- Is there a way to systematically search for the simulation parameters that produce the observed universe?
- What have we learned about galaxies at the end?
- How do we achieve the "Standard Model"?